

## CLAIMS

1 1. A proxy partition cache architecture for vending objects to requesting clients  
2 comprising:

3 a plurality of proxy cache servers for caching and vending objects from a storage  
4 location to a requesting client through an interconnected load-balancing mechanism  
5 adapted to selectively provide client requests to each of the plurality of proxy cache serv-  
6 ers based upon load-balancing considerations;

7 a network segment interconnecting each of the proxy cache servers so as to enable  
8 data to pass between the proxy cache servers; and

9 a mechanism in each server of the plurality of proxy cache servers adapted to (a)  
10 receive a request for an object from the load-balancing mechanism, (b) identify a discrete  
11 server of the plurality of proxy cache servers that is designated to contain the object,  
12 based upon an address of the object, (c) referring the request to the discrete server and (d)  
13 return the object from the discrete server to the server receiving the request from the  
14 load-balancing mechanism for vending to the client.

1 2. The proxy partition cache architecture as set forth in claim 1 wherein the mecha-  
2 nism includes means for the request to the discrete server to be referred over the network  
3 segment unless the discrete server and the receiving server are identical, whereby the re-  
4 quest is optimally processed on the receiving server.

5 3. The proxy partition cache architecture as set forth in claim 1 wherein the mecha-  
6 nism includes means for performing a hash function on the address of the object.

1 4. The proxy partition architecture as set forth in claim 3 wherein the mechanism  
2 includes a means for performing a modulo function on a hash of the address with respect  
3 to a number of proxy cache servers in the plurality of proxy cache servers.

1 5. The proxy partition architecture as set forth in claim 1 wherein the load-balancing  
2 mechanism comprises a network switch interconnected to the proxy cache servers.

1 6. The proxy partition architecture as set forth in claim 1 wherein the mechanism  
2 includes means for delivering the request to the discrete server and for tunneling the ob-  
3 ject through the server receiving the request from the discrete server to the load-balancing  
4 mechanism.

1 7. The proxy partition architecture as set forth in claim 1 wherein the mechanism  
2 includes means for enabling a redirection of the client request to the discrete server from  
3 the server receiving the request from the load-balancing mechanism.

1 8. The proxy partition architecture as set forth in claim 7 wherein the mechanism  
2 includes means for notifying the discrete server to receive the referred request so that a  
3 client redirection is expected, and allowing the discrete server to reject the request when  
4 the request is not recognized as a referred request by the mechanism.

1 9. The proxy partition architecture as set forth in claim 1 wherein the mechanism  
2 includes means for allowing vending of a file from a server other than the one selected by  
3 the modulo mechanism based upon a directory entry that indicates a current server loca-  
4 tion and data freshness of the file.

1 10. A method for address-partitioning a proxy cache defined by a group of proxy  
2 cache servers for vending objects to requesting clients comprising:  
3 caching and vending objects from a plurality of proxy cache servers to a request-  
4 ing client through an interconnected load-balancing mechanism, including selectively  
5 providing client requests to each of the plurality of proxy cache servers based upon load-  
6 balancing considerations;

7           interconnecting each of the proxy cache servers through a network segment so as  
8   to enable data to pass between the proxy cache servers; and

9           in each server of the plurality of proxy cache servers, (a) receiving a request for  
10   an object from the load-balancing mechanism, (b) identifying a discrete server of the plu-  
11   rality of proxy cache servers that is designated to contain the object, based upon an ad-  
12   dress of the object, (c) referring the request to the discrete server over the network seg-  
13   ment and (d) returning the object from the discrete server to the server receiving the re-  
14   quest from the load-balancing mechanism for vending to the client.

1   11.    The method as set forth in claim 10 further comprising referring the request to the  
2   discrete server over the network segment unless the discrete server and the receiving  
3   server are identical, whereby the request is optimally processed on the receiving server.

1   12.    The method as set forth in claim 10 wherein the step of identifying includes per-  
2   forming a hash function on the address of the object.

1   13.    The method as set forth in claim 12 wherein the step of identifying includes per-  
2   forming a modulo function on a hash of the address with respect to a number of proxy  
3   cache servers in the plurality of proxy cache servers.

1   14.    The method as set forth in claim 10 wherein the load-balancing mechanism com-  
2   prises a Layer 4 (L4) switch.

1   15.    The method as set forth in claim 10 wherein the steps of referring and returning  
2   each include delivering the request to the discrete server and tunneling the object through  
3   the server receiving the request from discrete server to the load-balancing mechanism.

1 16. The method as set forth in claim 10 further comprising allowing vending of a file  
2 from a server other than the one selected by the modulo mechanism based upon of a di-  
3 rectory entry that indicates a current server location and data freshness of the file.

1 17. The method as set forth in claim 10 wherein the steps of referring and returning  
2 include redirecting of the client request to the discrete server from the server receiving  
3 the request from the load-balancing mechanism.

1 18. The method as set forth in claim 17 further comprising notifying the discrete  
2 server to receive the referred request so that a client redirection is expected, and allowing  
3 the discrete server to reject the request as an unallowed external request when the request  
4 is not recognized as allowable.

1 19. A computer-readable medium including program instructions for address-  
2 partitioning a proxy cache defined by a group of proxy cache servers for vending objects  
3 to requesting clients, the computer-readable medium including instructions for perform-  
4 ing the steps of:

5 caching and vending objects from a plurality of proxy cache servers to a request-  
6 ing client through an interconnected load-balancing mechanism, including selectively  
7 providing client requests to each of the plurality of proxy cache servers based upon load-  
8 balancing considerations;

9 interconnecting each of the proxy cache servers through a network segment so as  
10 to enable data to pass between the proxy cache servers; and

11 in each server of the plurality of proxy cache servers, (a) receiving a request for  
12 an object from the load-balancing mechanism, (b) identifying a discrete server of the plu-  
13 rality of proxy cache servers that is designated to contain the object, based upon an ad-  
14 dress of the object, (c) referring the request to the discrete server over the network seg-

H:\112\024\0062C1\PROSECUT\PATAPP.doc 06/07/01 10:55 AM